

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-251412

(43)Date of publication of application : 06.09.2002

(51)Int.Cl.

G06F 17/30

G06F 12/00

(21)Application number : 2001-047027

(71)Applicant : CANON INC

(22)Date of filing : 22.02.2001

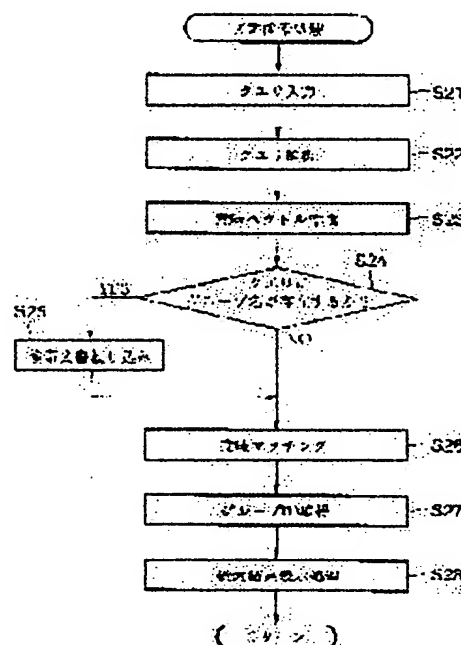
(72)Inventor : FUNAKOSHI MASANOBU

(54) DOCUMENT RETRIEVING DEVICE, METHOD, AND STORAGE MEDIUM

(57)Abstract:

PROBLEM TO BE SOLVED: To reduce processing time in document retrieving processing, where multi-dimensional vector matching is used, even when a large amount of document data is retrieved.

SOLUTION: Document data retrieving is conducted using a document database in which pairs, each composed of document data and a semantic vector representing its contents with a multi-dimensional vector, are stored. Document data stored in the document database is classified into groups according to similar semantic vectors, and the groups are stored in a group table. When a query described in a natural language statement is entered, a language analysis is made for the query to generate the semantic vector of the query (S21-S23). Based on the language analysis result, a group is selected from the group table (S24, S25). For document data belonging to the selected group, matching between the semantic vector stored in the selected group and the semantic vector of the query is performed to retrieve for a document and then the result is output (S26-S28).



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号
特開2002-251412
(P2002-251412A)

(43)公開日 平成14年9月6日(2002.9.6)

(51)Int.Cl. ⁷	識別記号	F I	テマコード*(参考)
G 0 6 F 17/30	3 5 0	G 0 6 F 17/30	3 5 0 C 5 B 0 7 5
	1 7 0		1 7 0 A 5 B 0 8 2
	3 3 0		3 3 0 C
12/00	5 2 0	12/00	5 2 0 E

審査請求 未請求 請求項の数16 O L (全 11 頁)

(21)出願番号 特願2001-47027(P2001-47027)

(22)出願日 平成13年2月22日(2001.2.22)

(71)出願人 000001007

キヤノン株式会社

東京都大田区下丸子3丁目30番2号

(72)発明者 船越 正伸

東京都大田区下丸子3丁目30番2号 キヤ
ノン株式会社内

(74)代理人 100076428

弁理士 大塚 康徳 (外1名)

Fターム(参考) 5B075 ND03 ND04 NK06 NK32 NK43

PP02 PP14 PP24 PQ36 PQ46

PR06 QM05 UU06

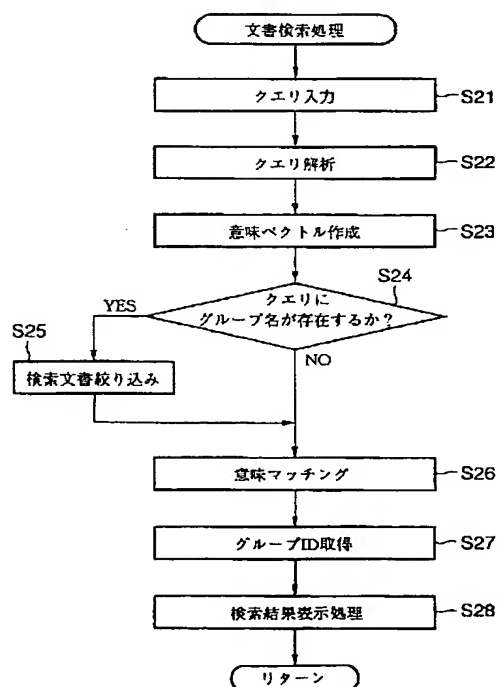
5B082 AA00 EA07

(54)【発明の名称】 文書検索装置および方法ならびに記憶媒体

(57)【要約】

【課題】多次元ベクトルマッチングを用いた文書検索処理において、膨大な文書データが検索対象となった場合でもその処理時間を短く保つ。

【解決手段】文書データと、その内容を多次元ベクトルで表現する意味ベクトルとを対にして登録した文書データベースを用いて文書データの検索を行う。ここで、文書データベースに登録された文書データは、その意味ベクトルの近いグループに分類されてグループテーブルに登録される。自然言語文でクエリが入力されると、これを言語解析して、クエリの意味ベクトルを生成する(S21～S23)。そして、この言語解析の結果に基づいて、グループテーブルに登録されているグループを選択する(S24, S25)。そして、選択されたグループに属する文書データについて、データベースに登録された意味ベクトルとクエリの意味ベクトルとのマッチングを行って文書を検索し、結果を出力する(S26～S28)。



【特許請求の範囲】

【請求項1】 文書データと、その内容を多次元ベクトルで表現する意味ベクトルとを対にして格納する格納手段と、
前記格納手段に格納された文書データを、その意味ベクトルに基づいてグループに分類して管理する管理手段と、
検索条件として入力された自然言語文を言語解析して、検索条件文の意味ベクトルを生成する生成手段と、
前記言語解析の結果に基づいて、前記管理手段で管理されているグループを選択する選択手段と、
前記選択手段で選択されたグループに属する文書データについて、前記格納手段に格納された意味ベクトルと前記検索条件文の意味ベクトルとのマッチングを行って文書を検索する検索手段とを備えることを特徴とする文書検索装置。

【請求項2】 前記管理手段は、
前記格納手段によって格納された各文書データの意味ベクトルを用いて自己組織化マップを生成し、
前記格納手段に格納された各文書データの意味ベクトルに基づいて、前記自己組織化マップ上に各文書データを配置し、
前記自己組織化マップ上で近隣に配置された文書データで1つのグループを形成させることを特徴とする請求項1に記載の文書検索装置。

【請求項3】 前記管理手段は、各グループ毎に当該グループを代表する単語によるグループ名を割り当て、
前記選択手段は、前記言語解析処理の結果得られた、前記検索文を含む単語と一致するグループ名を有するグループを選択することを特徴とする請求項1または2に記載の文書検索装置。

【請求項4】 前記管理手段は、
グループに属する文書データの意味ベクトルに基づいて、当該グループの意味ベクトルを設定し、
前記グループの意味ベクトルに基づいて取得された単語を前記グループ名とすることを特徴とする請求項3に記載の文書検索装置。

【請求項5】 前記管理手段は、
前記グループの意味ベクトルに基づいて複数の単語を取得し、取得された複数の単語のうち、シソーラスツリー上で中間部に位置する単語を前記グループ名として選択することを特徴とする請求項4に記載の文書検索装置。

【請求項6】 前記検索手段による検索結果として、前記マッチングの度合いが高い順に所定数の文書データを示す情報を出力する出力手段を更に備えることを特徴とする請求項1に記載の文書検索装置。

【請求項7】 前記出力手段は、前記文書データの各々に、当該文書データが属するグループを特定する情報を付加させて出力することを特徴とする請求項6に記載の文書検索装置。

【請求項8】 文書データと、その内容を多次元ベクトルで表現する意味ベクトルとを対にして格納する格納手段を用いた文書検索方法であって、
前記格納手段に格納された文書データを、その意味ベクトルに基づいてグループに分類して管理する管理工程と、
検索条件として入力された自然言語文を言語解析して、検索条件文の意味ベクトルを生成する生成工程と、
前記言語解析の結果に基づいて、前記管理工程で管理されているグループを選択する選択工程と、
前記選択工程で選択されたグループに属する文書データについて、前記格納手段に格納された意味ベクトルと前記検索条件文の意味ベクトルとのマッチングを行って文書を検索する検索工程とを備えることを特徴とする文書検索方法。

【請求項9】 前記管理工程は、
前記格納手段によって格納された各文書データの意味ベクトルを用いて自己組織化マップを生成し、
前記格納手段に格納された各文書データの意味ベクトルに基づいて、前記自己組織化マップ上に各文書データを配置し、
前記自己組織化マップ上で近隣に配置された文書データで1つのグループを形成させることを特徴とする請求項8に記載の文書検索方法。

【請求項10】 前記管理工程は、各グループ毎に当該グループを代表する単語によるグループ名を割り当て、
前記選択工程は、前記言語解析処理の結果得られた、前記検索文を含む単語と一致するグループ名を有するグループを選択することを特徴とする請求項8または9に記載の文書検索方法。

【請求項11】 前記管理工程は、
グループに属する文書データの意味ベクトルに基づいて、当該グループの意味ベクトルを設定し、
前記グループの意味ベクトルに基づいて取得された単語を前記グループ名とすることを特徴とする請求項10に記載の文書検索方法。

【請求項12】 前記管理工程は、
前記グループの意味ベクトルに基づいて複数の単語を取得し、取得された複数の単語のうち、シソーラスツリー上で中間部に位置する単語を前記グループ名として選択することを特徴とする請求項11に記載の文書検索方法。

【請求項13】 前記検索工程による検索結果として、前記マッチングの度合いが高い順に所定数の文書データを示す情報を出力する出力工程を更に備えることを特徴とする請求項8に記載の文書検索方法。

【請求項14】 前記出力工程は、前記文書データの各々に、当該文書データが属するグループを特定する情報を付加させて出力することを特徴とする請求項13に記載の文書検索方法。

【請求項15】 請求項8乃至14のいずれかに記載の文書検索方法をコンピュータによって実現させるための制御プログラム。

【請求項16】 請求項8乃至14のいずれかに記載の文書検索方法をコンピュータによって実現させるための制御プログラムを格納する記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、言語解析による文書検索装置及び方法ならびにこれを記憶した媒体に関する。

【0002】

【従来の技術】インターネットの爆発的な普及と共に、その入り口となるパーソナルコンピュータ（PC）や携帯情報機器（PDA）の数も増大し、電子メールなどを利用した電子化文書によるコミュニケーションは我々の日常生活に浸透しつつある。このような状況において電子化文書の数は現在もなお増え続ける一方である。このため、膨大な文書データの中から必要な文書を迅速に検索する技術はますます重要視されており、より使いやすいインターフェースを持つ文書検索システムが登場している。

【0003】この種の文書検索システムに、例えば、意味概念を用いた文書検索システムがある。これは、文書データベースにおいて、各文書データとその内容の意味概念を示す多次元ベクトルとを一对一に対応付けて格納しておき、ユーザが日常使っている自然言語を用いた簡単な質問文が入力されると、この質問文を言語解析してその意味概念を表す多次元ベクトルに変換し、文書データベース中の多次元ベクトルとマッチングをとることによって検索するものである。この手法によれば、ユーザが探している文書を自然言語を用いたインターフェースで検索することができる。

【0004】

【発明が解決しようとする課題】しかしながら、多次元ベクトルのマッチングは、キーワードマッチングなどと比較すると計算量が大きくなるため、検索対象となる文書データの数が膨大になると、検索処理に多大な時間がかかってしまうという課題があった。

【0005】また、多次元ベクトルのマッチングによる検索では、場合によっては検索結果の中に全く関係のない文書が混じってしまい、検索結果の精度が上がらないという課題もあった。

【0006】本発明は、上記の課題に鑑みてなされたものであり、多次元ベクトルマッチングを用いた文書検索処理において、膨大な文書データが検索対象となった場合でもその処理時間を短く保つことを目的とする。

【0007】

【課題を解決するための手段】上記の課題を達成するための本発明による文書検索装置は以下の構成を備える。

すなわち、文書データと、その内容を多次元ベクトルで表現する意味ベクトルとを対にして格納する格納手段と、前記格納手段に格納された文書データを、その意味ベクトルに基づいてグループに分類して管理する管理手段と、検索条件として入力された自然言語文を言語解析して、検索条件文の意味ベクトルを生成する生成手段と、前記言語解析の結果に基づいて、前記管理手段で管理されているグループを選択する選択手段と、前記選択手段で選択されたグループに属する文書データについて、前記格納手段に格納された意味ベクトルと前記検索条件文の意味ベクトルとのマッチングを行って文書を検索する検索手段とを備える。

【0008】また、上記の課題を達成するための本発明による文書検索方法は、文書データと、その内容を多次元ベクトルで表現する意味ベクトルとを対にして格納する格納手段を用いた文書検索方法であって、前記格納手段に格納された文書データを、その意味ベクトルに基づいてグループに分類して管理する管理工程と、検索条件として入力された自然言語文を言語解析して、検索条件文の意味ベクトルを生成する生成工程と、前記言語解析の結果に基づいて、前記管理工程で管理されているグループを選択する選択工程と、前記選択工程で選択されたグループに属する文書データについて、前記格納手段に格納された意味ベクトルと前記検索条件文の意味ベクトルとのマッチングを行って文書を検索する検索工程とを備える。

【0009】

【発明の実施の形態】以下、添付の図面を参照しながら本発明の好適な実施形態を説明する。

【0010】本実施形態では、文書データを登録する際に、文書の意味内容を示す多次元ベクトルを作成し、この多次元ベクトルを利用して予め文書を意味概念によって自動分類して管理する。そして、検索時のユーザクエリにこの分類と合致する言葉が現れた場合に、まず検索対象をこの分類に絞り込み、多次元ベクトルマッチングを行う。この検索対象の絞り込みにより検索結果の精度を高めると同時に、文書検索時の計算量を減らし、正確で迅速な文書検索を実行する。以下、本実施形態について詳細に説明する。

【0011】図1は本実施形態による文書検索処理を実行するコンピュータ装置100の構成を示すブロック図である。図1の構成において、CPU101はマイクロプロセッサであり、文書検索処理のための演算、論理判断等を行い、PCIバス102を介して接続された各構成要素を制御する。PCIバス102はCPU101の制御の対象とする構成要素を指示するアドレス信号を転送し、CPU101の制御の対象とする各構成要素のコントロール信号を転送し、各構成機器相互間のデータ転送を行う。

【0012】ROM103は読み出し専用の固定メモリで

ある。ROM103には本実施形態の構成における基本I/Oプログラムが格納される。また、RAM104は、書込み可能のランダムアクセスメモリであって、各構成要素からの各種データの一時記憶と、本実施形態における各種処理が記述されたプログラムが格納され、このプログラムに基づいてCPU101が各種処理を行う。

【0013】DVDD105はDVDドライブである。DVDメディア(DVD-MEDIA)106に記録されているプログラムやデータはこのDVDドライブ105を通じて本システムにロードされる。また、DISK107に蓄えられた各種データをDVDドライブ105を通じてDVD-MEDIA106に書き込むことができる。なお、DVD-MEDIA106は、具体的にはDVD-ROM、DVD-RAM、DVD-R、DVD-RW、DVD-VIDEO、DVD-AUDIOなどのDVD規格のメディアを総称したものである。本実施形態において、DVD-MEDIA106は文書データやその関連データ、もしくはプログラムなどの大容量データの読み書きに用いられる。

【0014】INPUTC108は入力コントローラである。キーボード(KB)109やポインティングデバイス(PD)110から送られてくる入力信号は、このコントローラによって適宜適切な信号に変換された後、PCIバス102を経由してCPU101に送信される。

【0015】KB109は、アルファベットキー、ひらがなキー、カタカナキー等の文字記号入力キー、及び、カーソル移動を指示するカーソル移動キー等のような各種のファンクションキーを備えている。PD110は、マウスやトラックボールなどのポインティングデバイスであり、表示画面上のカーソルやボタンなどを指摘するために使用される。

【0016】DISK107はデータやプログラム等を記憶するための外部メモリである。データやプログラム等は必要に応じて保管され、また、保管されたデータやプログラムはキーボードの指示により、必要な時に呼び出される。本実施形態における文書データベースは主にこのDISK107上に実装される。

【0017】VIDEO111はビデオコントローラである。PCIバス102を介して表示用のデータがここに蓄えられるとともに、表示用の信号に変換されて表示装置(DISPLAY)112に出力される。DISPLAY112には、陰極線管や液晶などが用いられ、各種処理の結果や装置の状態、ユーザに対するメッセージなどを表示する。

【0018】DEVC113はデバイスコントローラである。PCIバス102を介して伝達されるCPU101の指示によって、このコントローラに接続されている機器を制御し、また、接続されている機器が出力する信

号やデータをPCIバス102を介してCPU101やDISK107に適宜伝達する。SCAN114はスキャナである。DEVC113からの指示によって、光学的な方法によってここにセットされた原稿をスキャンし、原稿画像を読み取り、これをDEVC113に出力する。

【0019】NI115はネットワークインターフェースであり、本実施形態の文書検索システムをLANやインターネット116などを経由して外部のシステムと接続するための機器である。本実施形態の文書検索システムは、この接続を経由して、信号やデータを外部のシステムと送受信することが可能である。

【0020】MIX117はミキサーである。PCIバス102を介して音声出力用のデータがここに送られると、MIX117はこれらの信号を合成しかつ音声出力用の信号に変換してスピーカ(SPK)118に出力する。SPK118は、処理結果や装置の状態、ユーザに対するメッセージ、音楽などを音声で出力する。

【0021】かかる各構成要素からなる本実施形態のコンピュータ装置100においては、キーボード109やポインティングデバイス110からの各種の入力に応じた各種処理を実行させる。すなわち、キーボード109やポインティングデバイス110から入力信号が供給されると、INPUTC108を経由して、インタラプト信号がCPU101に送られ、CPU101がROM103内に記憶してある各種の制御信号を読み出し、それらの制御信号に従って、各種の制御が行われる。

【0022】本実施形態では、コンピュータが基本I/Oプログラム、OS、及び本文書検索処理プログラムをCPU101が実行することによって、文書検索装置として動作する。基本I/OプログラムはROM103中に書き込まれており、OSはDISK107に書き込まれている。そして、本システムの電源がONにされると、基本I/Oプログラム中のIPL(イニシャルプログラムローディング)機能によりDISK107からOSがRAM104に読み込まれ、CPU101によるOSの動作が開始される。なお、文書検索処理プログラムは、図9～図15のフローチャートによって示される文書検索処理手順をCPU101によって実現させるためのプログラムコードである。

【0023】図2は、本文書検索処理プログラム及び関連データをDVD-MEDIA106に記録したときのDVD-MEDIA106の内容の構成図である。本実施形態において、文書検索処理プログラム及び関連データはDVD-MEDIA106に記録されている。図示のようにDVD-MEDIA106の先頭領域には、このDVD-MEDIAのボリューム情報201とディレクトリ情報202が記録されており、その後このDVD-MEDIA106のコンテンツである本実施形態の文書検索処理プログラム(実行ファイル)203と、文

書検索処理関連データ204が記録されている。

【0024】図3はコンピュータ装置100と本文書検索処理プログラムが記録されたDVD-MEDIA106の模式図である。DVD-MEDIA106に記録された文書検索処理プログラム203および関連データ204は、図3に示したようにDVD-MEDIAドライブ(DVDD)105を通じて本システムにロードすることができる。すなわち、このDVD-MEDIA106をDVDD105にセットすると、OS及び基本I/Oプログラムの制御のもとに本文書検索処理プログラムおよび関連データがDVD-MEDIA106から読み出され、RAM104にロードされて動作可能となる。

【0025】図4は、本文書検索処理プログラムがRAM104にロードされ実行可能となった状態のメモリマップを示す。コンピュータ装置100の立ち上げ時にROM103やDISK107よりロードされる基本I/Oプログラム411とOS412が格納されている。また、DVD-MEDIA106からロードされた文書検索処理プログラム203とその関連データ204がそれぞれ文書検索処理プログラム413、関連データ414として格納される。また、ワークエリア415には、意味概念辞書MDIC421、検索結果バッファRBUF422、結果出力数COUNT424、グループテーブルGTBL423が存在している。

【0026】図5は、本実施形態における意味概念辞書MDICの構成例を説明した図である。

【0027】本実施形態における意味概念辞書(MDIC)421には、図示したように、単語ID501と、単語表記502と、その意味表現である多次元ベクトル503のリストで構成される。図5において、表の行が単語一個分のデータに相当する。このうち、単語ID501は、本実施形態の文書検索システムにおいて、各単語を識別管理するために、各単語に対して一意に割り振られている番号である。また、単語表記502は、各単語の表記を表す文字列である。また、意味ベクトル503は、各単語が持つ意味概念を多次元ベクトルで表現したものであり、予め定められている。

【0028】なお、本実施形態における意味ベクトルの各次元は、名詞のシソーラスなどを参照して上位概念に相当する単語を適宜抽出して設定する。また、各単語の意味ベクトルは国語辞書などの語義文を言語解析して、各次元として選択された概念の合成に落とし込むことで作成するが、このような手法は意味処理として一般的であり、公知であるので、ここでは詳述しない。

【0029】意味概念辞書(MDIC)421は以上の構成によって、単語ID501もしくは単語表記502をキーにして検索され、格納されている各情報を参照することが可能である。本実施形態において、意味概念辞書(MDIC)421は、後述する各処理において適宜検索、参照される。

【0030】図6は、本実施形態における検索結果バッファ(RBUF)の構成例を説明した図である。本実施形態における検索結果バッファ(RBUF)422には、文書検索処理において検索されたデータのうち、意味マッチングの度合いが高いデータから順に結果出力数COUNT424の数だけ格納される。図5は、COUNT=10の場合の検索結果バッファ(RBUF)422の構成例である。

【0031】図6に示したように、一つの検索結果データは、順位601と、意味近似度602と、意味ベクトル603と、文書ID604と、グループID605によって構成される。このうち、順位601は、RBUF422に格納されたデータの中における意味近似度が大きい順番である。また、意味近似度602は、各文書データに対応付けられた意味ベクトル603と、クエリを言語解析して作成した意味ベクトルとのマッチングを取ったときの度合いであり、本実施形態においてはパーセントで格納される。なお、2つの意味ベクトルのマッチングの度合いは、主にこれら2つのベクトルのなす角を算出することによって決定される。

【0032】意味ベクトル603は、各文書データ中のテキストを言語解析して意味情報を抽出し、多次元ベクトルで表現したものであり、登録された各文書について予めデータベース中に格納されている。この意味ベクトル603は、通常、テキストを形態素解析を用いて品詞分解し、普通名詞や固有名詞を抜き出し、これらの持つ意味ベクトルを意味概念辞書(MDIC)502から検索して取得後、構文情報によって重み付けして合成することによって作成される。なお、このような手法は自然言語処理分野において公知であるので、ここでは詳述しない。

【0033】また、文書ID604は、本実施形態の文書データベースに格納されている文書のIDであり、このIDを利用して実際の文書データを読み出すことが可能である。また、グループID605は、各文書が属する意味近似グループのIDであり、これを用いてグループテーブル(GTBL)423を参照することにより、この文書のグループ情報を取得することができる。なお、意味近似グループについては図7を参照して以下に説明する。

【0034】図7は、本実施形態におけるグループテーブル(GTBL)の構成例を説明した図である。図7に示したように、本実施形態におけるグループテーブル(GTBL)423は、グループID701と、グループ名702と、データ数703と、文書IDリスト704と、グループの意味ベクトル705によって構成される。このうち、グループID701は、本実施形態における文書データベースに登録してある文書データを、後述する手法によって意味的に分類することによって作成されるグループに一意に割り振られる番号である。な

お、このグループIDによって特定される1つのグループが1つの意味近似グループを構成する。このグループID701はRBUF422のグループID605に対応している。また、グループ名702は、本実施形態における各意味近似文書グループに与えられたユニークな名前である。このグループ名を決定する処理は図15を用いて後述する。また、データ数703は、このグループに属している文書データの数である。

【0035】また、文書IDリスト704は、このグループに属する文書データのIDのリストである。グループの意味ベクトル705は、当該グループに属する文書データの全ての意味ベクトルとのなす角度が最小となる多次元ベクトルである。グループテーブル(GTBL)423は、後述する意味近似グループ作成処理によって作成され、主に文書検索時の検索対象データの絞り込みに使用される。図7は、PC関連のニュース記事を自動分類した場合のグループテーブルGTBLの一例を示している。

【0036】図8は、本実施形態における文書データベースに格納されている一つのデータ構造の構成例を説明した図である。図8に示したように、本実施形態における文書データベース中の1つのデータは、文書ID801、文書データポインタ802、キーワード803、意味ベクトル804によって構成される。

【0037】ここで、文書ID801はこのデータ構造そのもののインデックス番号であると同時に、データベース中の文書データ自体のインデックス番号である。文書ID801はRBUF422における文書ID604、GTBL702における文書IDリストで用いられる文書IDである。また、文書データポインタ802は、実際の文書データを指すポインタであり、このポインタを利用することにより実際の文書データにアクセスすることができる。また、キーワード803は、文書データ中に現れる代表的な単語のリストである。また、意味ベクトル804は文書データの意味を多次元ベクトルで表現したもので、RBUF422に格納される意味ベクトル603と同じものである。このようなデータ構造を取ることによって、実際の文書データにアクセスすることなく検索を高速に行うことが可能となる。

【0038】以上の構成を備えた本実施形態の文書検索装置の動作について、以下、詳細に説明する。

【0039】図9は、本実施形態における文書検索処理の全体を説明するフローチャートである。ステップS1は、ユーザのKB109もしくはPD110の操作を受け付け、これをコマンドとして解釈するユーザコマンド入力処理である。このような処理はPCなどを利用したシステムにおいて一般的なユーザインターフェース処理であり、公知であるのでここでは詳述しない。処理を終えると、ステップS2へ進む。

【0040】ステップS2は、ステップS1で解釈した

コマンドを判定して、各種処理に分岐するユーザコマンド判定処理である。ユーザコマンドが文書登録を指示している場合は、ステップS3へ進む。また、ユーザコマンドが文書検索を指示している場合は、ステップS4へ進む。また、ユーザコマンドがその他の処理を指示している場合は、ステップS5へ進む。

【0041】ステップS3では、検索対象となる文書データベースに新たに文書データを登録する文書登録処理を行う。この処理は、図10のフローチャートを用いて後述する。文書登録処理を終えると、ステップS1へ戻る。ステップS4では、ユーザクエリを受け付けて、そのクエリに合致するデータベース中の文書データを検索する文書検索処理を実行する。この文書検索処理は、図11のフローチャートを用いて後述する。処理を終えると、ステップS1へ戻る。ステップS5では、表示のカスタマイズ、システムの各種設定変更、ユーザ情報登録などの、その他の処理を行う。これらの処理はこのようなシステムにおいて一般的であり、公知であるので詳述しない。処理を終えると、ステップS1へ戻る。

【0042】なお、本実施形態のシステムにおいて、ステップS1へ処理が進むとき、新規登録文書数がある一定値を超えた、グループ情報の更新指示がされたなどの、諸処の条件が満たされた場合、システムに割り込みがかかり、ステップS6へ処理が進む。ステップS6は、その時点で本実施形態の文書データベースに登録されている文書をその意味によって新たにグループ分けして、図7に示したようなグループテーブルGTBLを作成する、意味近似文書グループ作成処理を行う。なお、この処理の詳細は、図12のフローチャートを用いて後述する。意味近似文書グループ作成処理を終えると、ステップS1へ処理が戻り、割り込みから復帰する。

【0043】図10は、図9のステップS3における文書登録処理を詳細化したフローチャートである。本処理により、図8に示す形態で文書が文書データベースに登録されることになる。

【0044】ステップS11では、ユーザが登録する文書を指定する登録文書指定処理を実行する。この処理は、実際にはGUIを利用した文書ファイルのドラッグ・ドロップであったり、あるいは、登録ファイル名のリストを格納したテキストファイルの指定でも良い。この種の処理は本実施形態のようなシステムのユーザインターフェースとして極めて一般的であるので、詳述しない。処理を終えると、ステップS12へ進む。

【0045】ステップS12では、ステップS11で指定された文書データのうち、テキストの主要部分を抽出して言語解析する主要文解析処理が行われる。ここで、テキスト主要部分とは、テキストの大意が記述されている可能性が高い部分を抽出したものであり、文書の論理、レイアウト構造を利用して経験的に抽出される。主要文解析処理を終えると、ステップS13へ進む。

【0046】ステップS13では、ステップS12で得られた言語解析（主要文解析処理）の結果を利用して、この文書において特徴的に現れるキーワードを、重み付けして抽出するキーワード抽出処理が行われる。このようなキーワード抽出処理は言語処理の分野で一般的に行われており、公知であるので詳述しない。キーワード抽出処理を終えると、ステップS14へ進む。

【0047】ステップS14では、ステップS13で得られたキーワード情報を元に、この文書データの意味概念を表す多次元ベクトルを作成する文書意味ベクトル作成処理が行われる。本実施形態において、文書意味ベクトルの作成は、以下のように行われる。まず、ステップS13で得られた各キーワードを意味概念辞書MDIC421（図5）から検索し、各キーワードの意味ベクトルを取得する。次に、各意味ベクトルにステップS13で得られた重み付けを施した後、全ての意味ベクトルを合成する。処理を終えると、ステップS15へ進む。

【0048】ステップS15では、ステップS13で得られたキーワード、ステップS14で得られた意味ベクトルの情報から図8に示した文書データを作成して文書データベースに登録する。なお、文書IDは、重複が発生しないように、当該システムによって自動的に割り振られる。また、同時に、実際の文書データ自体もデータベース上の別領域に登録され、このときに文書データポインタが得られ、データベースに登録される。データ登録処理を終えると、ステップS16へ進む。

【0049】ステップS16では、ステップS14で作成した意味ベクトルを利用して、登録する文書が属するべき意味近似文書グループを決定する。意味近似文書グループの決定は、登録する文書の有する意味ベクトルと、グループテーブル（GTBL）中の各グループの意味ベクトル705とのマッチングを取り、最も近似しているグループを決定することによって行われる。文書の属するグループが決定されると、図7で示したグループテーブル（GTBL）中の決定されたグループの文書IDリストに、ステップS15で割り振られた文書IDが格納される。処理を終えると、ステップS17へ進む。

【0050】ステップS17は、これまでの登録処理の結果をユーザに通知する登録結果通知処理である。例えば、これは結果を記載したダイアログボックスの表示であったり警告音の出力であるが、この種の処理はユーザインターフェースとして一般的であり、公知であるので詳述しない。処理を終えると、文書登録処理を終了する。

【0051】図11は、図9に示したステップS4の文書検索処理の詳細を示すフローチャートである。

【0052】ステップS21では、文書を検索するためのクエリ（質問文）をユーザに入力してもらうクエリ入力処理を実行する。クエリは、ユーザが探している文書の内容を表現する簡単な文であり、例えば、「PCの価

格下落について」のように自然言語にて入力する。このような処理は一般の検索処理において公知であり、ここでは詳細な説明は行わない。クエリ入力処理を終えると、ステップS22に進む。

【0053】ステップS22では、ステップS21で入力されたクエリを言語解析するクエリ解析処理が実行される。この処理では、ステップS21で入力されたクエリに、形態素解析、構文解析等の言語処理を行い、後の処理で利用しやすい形式である言語情報に変換する。なお、このような処理は一般の言語処理において公知であり、ここでは詳細な説明は行わない。クエリ解析処理を終えると、ステップS23に進む。

【0054】ステップS23では、ステップS22の結果得られた言語情報を利用して、クエリの意味ベクトルを作成する意味ベクトル作成処理が行われる。この意味ベクトル作成処理では、文書の意味ベクトルを作成する場合（ステップS14）と同様に、クエリ中に現れる単語の意味ベクトルを意味概念辞書（MDIC）421から検索して取得し、これらを構文情報によって重み付けしてから合成することで作成する。意味ベクトル作成処理を終えると、ステップS24へ進む。

【0055】ステップS24は、ステップS22で得られた言語情報を利用して、クエリ中にグループテーブル（GTBL）423に格納されている意味近似グループのグループ名が存在しているかどうか調べる処理である。この処理の結果、クエリ中に1つ以上のグループ名が存在する場合はステップS25へ進む。グループ名が1つも存在しない場合はステップS26へ進む。

【0056】ステップS25では、ステップS24の結果、クエリ中に存在するグループ名を持つ意味近似グループに属する文書のみに検索対象となる文書を絞り込む。検索対象をこの時点で絞り込むことによって、後の意味マッチング処理における処理量（計算量）を減らし、検索を高速化することができる。検索文書の絞り込み処理が終わるとステップS26へ進む。

【0057】ステップS26では、ステップS23において得られたクエリの意味ベクトルと、文書データベース中の意味ベクトルとのマッチングを行って、ユーザが探している文書を検索する意味マッチング処理が実行される。意味ベクトルのマッチングは、前述したように2つの多次元ベクトルのなす角を算出することによって行われる。なお、ステップS25によって検索対象となる文書が絞り込まれた場合は、検索対象となっている文書の意味ベクトルとのみマッチングが行われる。この処理の結果として、マッチング結果バッファ（RBUF）422には、データベース中のデータの中で、マッチングの度合いが上位のものから順に、結果出力数COUNTの数だけ検索結果データが格納される。なお、このような処理は一般のデータベース検索処理において公知であり、ここでは詳細な説明は行わない。以上の意味マッ

ング処理を終えると、処理はステップS27に進む。

【0058】ステップS27では、ステップS26によって出力された検索結果バッファ(RBUF)中の各文書データが属するグループのIDを、グループテーブル(GTBL)423を検索することによって取得する、グループID取得処理を実行する。この処理によって取得したグループIDは、検索結果バッファ(RBUF)422中のグループID605として格納される。処理を終えると、ステップS28へ進む。

【0059】そして、ステップS28では、ステップS27によって作成された検索結果バッファRBUF中の検索結果データに基づいて、ディスプレイに文書データを出力する検索結果表示処理である。この種の処理は表示を行うシステムにおいて一般的に行われており、公知であるので詳述しない。処理を終えると、文書検索処理を終了する。

【0060】図12は、図9のステップS6に示した意味近似文書グループ作成処理の詳細を示すフローチャートである。この処理は、各文書の多次元ベクトルを利用して予め文書を意味概念によって自動分類して管理するものであり、結果として図7に示したグループテーブル(GTBL)423が自動生成される。

【0061】ステップS31では、本実施形態における文書データベースに格納されている意味ベクトルによって自己組織化マップ学習を行う。この自己組織化マップ学習処理は、多次元ベクトル表現の集合を2次元マップに変換して分類する手法として一般的に行われており、公知であるが、ここではその概要を図13を用いて簡単に説明する。

【0062】図13は自己組織化マップの学習の例を示す図である。図13において、1はユニットU、2はUの近傍ユニットUcである。図13に示したように、自己組織化マップは正方形であり、 $N \times N$ 個(N は正の整数)のユニットで構成される。各ユニットは、意味ベクトルと同次元のパターンベクトルを持つ。

【0063】自己組織化マップの学習は、以下の手順で行う。まず、マップ中のユニットが持つ全てのパターンベクトルをゼロベクトルに初期化する。次に、文書データベース中の各意味ベクトルに対して、なす角が最小であるパターンベクトルを持つユニットUを決定する。次に、Uとその近傍ユニットUcが持つパターンベクトルを意味ベクトルに近づける。このとき、各ユニットが持つパターンベクトルに対して、ユニットUとの距離が近いほど、パターンベクトルを意味ベクトルに近づける度合いを高くする。初期化以外の以上の操作を予め決められた回数繰り返すと、2次元マップ上の各ユニットが持つベクトルパターンはある定値に収束し、自己組織化マップ学習が終了する。自己組織化マップ学習処理を終えると、ステップS32へ進む。

【0064】ステップS32では、ステップS31で学

習が済んだ自己組織化マップ上のユニットに文書データベースに登録されている各文書の意味ベクトルを配置する処理が行われる。各文書の意味ベクトルは、マップ学習時と同様に、パターンベクトルとのなす角が最小であるユニットに配置される。処理を終えると、ステップS33へ進む。

【0065】ステップS33では、ステップS32によって2次元マップ上に配置された各文書データを8連結手法によってグループ化し、グループを決定する処理が行われる。ここで、8連結手法とは、2次元マップ上で縦横斜めに隣接するユニット同士を同じグループとしてまとめる手法である。図14は、2次元マップ上において8連結手法によってグループを決定する例を示す図である。図14において、3はグループID=1のグループ、4はグループID=2のグループ、5はグループID=3のグループ、6はグループID=4のグループである。なお、マップ上の数字は配置された文書データの数を示す。処理を終えると、ステップS34へ進む。ステップS34では、ステップS43で決定したグループに属する文書データの全ての意味ベクトルとのなす角度が最小となる多次元ベクトルを求め、これをグループの意味ベクトルとする。

【0066】ステップS35では、ステップS34で決定された意味近似文書グループのグループ名を決定する、グループ名決定処理が行われる。この処理は、図15を用いて後述する。グループ名決定処理を終えるとステップS36へ進む。ステップS36は、ステップS33～ステップS35で作成されたグループの情報を参照して、グループテーブル(GTBL)423を作成する処理である。処理を終えると、意味近似文書グループ作成処理を終了する。なお、本例では、グループテーブル(GTBL)423において、文書データの総数の多い順に、グループIDを1から順に割り当てる。例えば、図14において、グループ3のマッチング数は $1+5+2+8+9=25$ で最も多く、図7に示すように、ID=1が割り当てられている。なお、グループ作成処理の際には、図6の検索結果は破棄されているものとする。

【0067】図15は、ステップS34におけるグループ名決定処理の詳細を説明するフローチャートである。図15において、ステップS41では、グループ名を決定するグループの意味ベクトル705と、意味概念辞書(MDIC)中の意味ベクトルとのマッチングを取り、ある閾値以上の類似度を持つ単語を当該グループのグループ名候補として複数個検索する。また、ここで用いる閾値は、予め本実施形態のシステムにおいて定まっているものとする。

【0068】次に、ステップS42において、ステップS41で選択したグループ名候補単語が、シソーラスツリー上においてどの位置にあるかを判定し、シソーラスツリー上の中間部に位置する単語をグループ名として選

択する。これは、シソーラスツリー上の上位（根の部分）に位置する単語をグループ名として選択すると、意味概念が大きくなりすぎるためそのグループの持つ概念がぼやけてしまうのでそれを避けるためである。また、逆に、シソーラスツリー上の下位（子葉）部分に位置する単語をグループ名として選択すると、その単語が持つ意味概念が狭すぎるために、グループ名にふさわしくない文書データがグループに紛れ込む可能性が高くなるので、子葉部分に位置する単語は候補から外すためである。

【0069】この処理において、シソーラスツリー上の中間部に位置する単語のうち、グループの意味ベクトル705との類似度が最も高い意味ベクトルを持つ単語がグループ名として決定される。処理を終えると、グループ名決定処理を終了する。

【0070】以上のように、本実施形態によれば、文書データを登録する際に、文書の意味内容を示す多次元ベクトルが作成され、このベクトルを利用して予め文書が意味概念によって自動分類され、意味概念の類似している文書グループを表現するのにふさわしい単語がグループ名として付加される。そして、検索時のユーザクエリにこの単語が現れた場合に、検索結果候補をこの単語をグループ名とする意味近似グループに絞り込むことによって、意味ベクトルのマッチング処理量を減らし、迅速な検索を実行することが可能となる。また、特に、グループ名としてシソーラスツリーの中間に位置する単語を選択することにより、検索精度を高めることが可能となる。

【0071】なお、本発明は上述した実施形態に限定されるものではない。例えば、上述の実施形態では、文書検索装置のバスとしてP C Iバスを採用しているが、I S AバスやV Lバスなどでもまったく同様な文書検索装置を構成することが可能である。また、上述の実施形態では、OSはD I S Kに格納されているが、OSをR O M上に格納しても同様な処理を行うことが可能である。

【0072】また、上述の実施形態では、DVD-M E D I Aから文書検索処理プログラムおよび関連データを直接R A Mにロードして実行させる例を示したが、このほかにD C D-M E D I Aから文書から文書検索処理プログラムおよび関連データを一旦D I S Kに格納（インストール）しておき、本文書検索処理プログラムを動作させるときにD I S KからR A Mにロードするようにすることも可能である。

【0073】また、上述の実施形態では、本文書検索処理プログラムを記憶する媒体としてDVD-M E D I Aを用いているが、それ以外にC D-M E D I A, M O, D F, I Cメモ리카ード、光磁気カードなどを用いても良い。更に文書検索処理プログラムをR O Mに記憶しておき、これをメモリマップの一部となすように構成し、直接C P Uで実行することも可能である。

【0074】また、本発明の目的は、前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記憶媒体（または記録媒体）を、システムあるいは装置に供給し、そのシステムあるいは装置のコンピュータ（またはC P UやM P U）が記憶媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。この場合、記憶媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記憶した記憶媒体は本発明を構成することになる。また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているオペレーティングシステム（O S）などが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

【0075】さらに、記憶媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張カードやコンピュータに接続された機能拡張ユニットに備わるメモリに書込まれた後、そのプログラムコードの指示に基づき、その機能拡張カードや機能拡張ユニットに備わるC P Uなどが実際の処理の一部または全部を行い、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

【0076】その他、本発明はその要旨を逸脱しない範囲で種々変形して実施することができる。

【0077】

【発明の効果】以上説明したように、本発明によれば、多次元ベクトルマッチングを用いた文書検索処理において、膨大な文書データが検索対象となった場合でもその処理時間を短く保つことができる。

【図面の簡単な説明】

【図1】本実施形態による文書検索処理を実行するコンピュータ装置100の構成を示すブロック図である。

【図2】本文書検索処理プログラム及び関連データをDVD-M E D I A 106に記録したときのDVD-M E D I A 106の内容の構成図である。

【図3】コンピュータ装置100と本文書検索処理プログラムが記録されたDVD-M E D I A 106の模式図である。

【図4】本文書検索処理プログラムがR A M 104にロードされ実行可能となった状態のメモリマップを示す図である。

【図5】C O U N T = 10の場合の検索結果バッファ（R B U F）422の構成例を説明した図である。

【図6】本実施形態における検索結果バッファ（R B U F）の構成例を説明した図である。

【図7】本実施形態におけるグループテーブル（G T B L）の構成例を説明した図である。

【図8】本実施形態における文書データベースに格納されている一つのデータ構造の構成例を説明した図である。

【図9】本実施形態における文書検索処理の全体を説明するフローチャートである。

【図10】図9のステップS3における文書登録処理を詳細化したフローチャートである。

【図11】図9に示したステップS4の文書検索処理の詳細を示すフローチャートである。

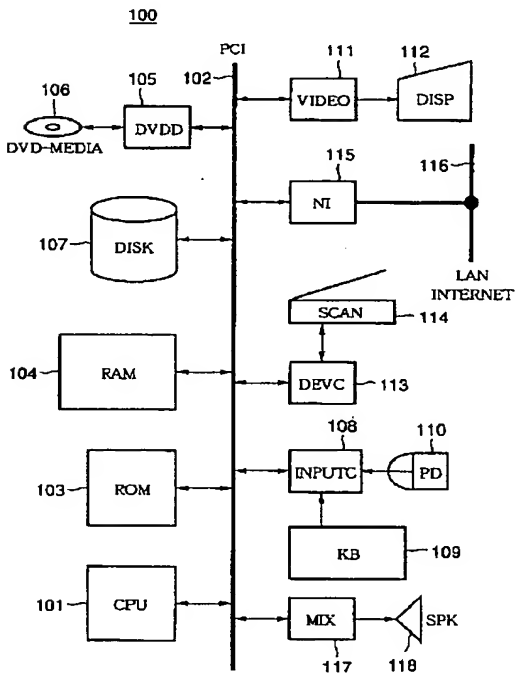
【図12】図9のステップS6に示した意味近似文書グループ作成処理の詳細を示すフローチャートである。

【図13】自己組織化マップの学習の例を示す図である。

【図14】自己組織化マップの学習において、2次元マップ上で8連結手法によってグループを決定する例を示す図である。

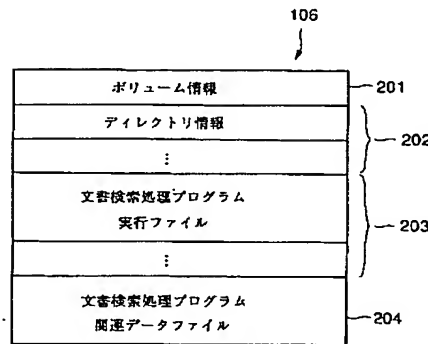
【図15】ステップS34におけるグループ名決定処理の詳細を説明するフローチャートである。

【図1】

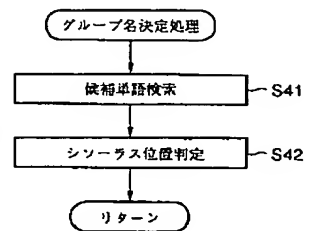


【図3】

【図2】



【図15】

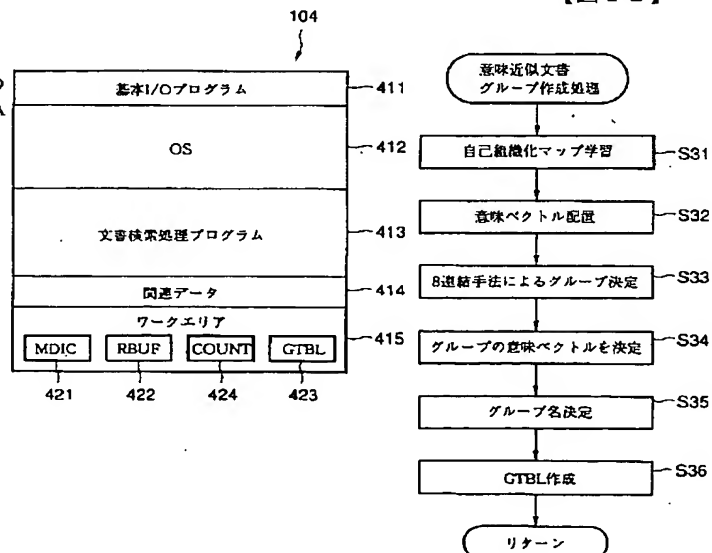


【図5】

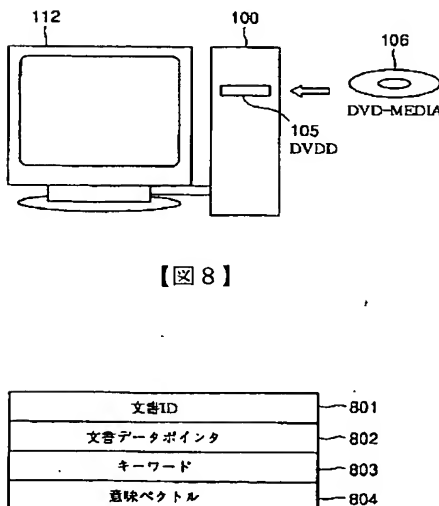
単語ID	単語表記	意味ベクトル
1	愛	(5,8,2,0,...)
2	挨拶	(6,1,0,0,...)
...
...

【図4】

【図12】



【図8】



【図6】

422

順位	意味近似度	意味ベクトル	文書ID	グループID
1	80	(1,0,4,8,...)	21365	1
2	75	(2,0,0,3,...)	90235	1
:	:	:	:	:
10	35	(0,0,0,5,...)	38106	3

RBUF 601 602 603 604 605

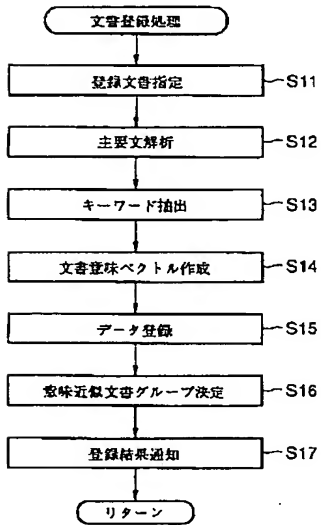
【図7】

423

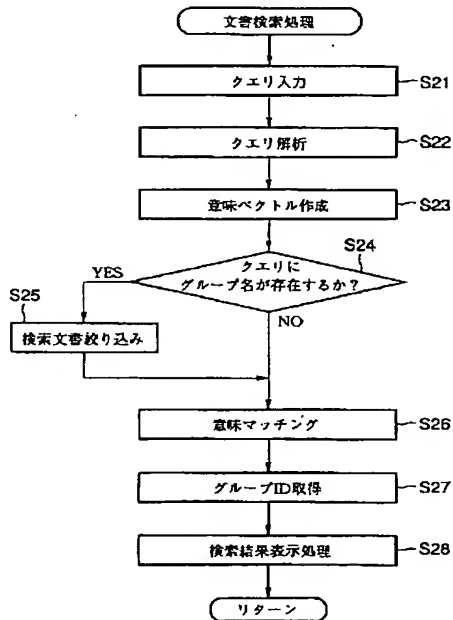
グループID	グループ名	データ数	文書IDリスト	意味ベクトル
1	インターネット	52	21365, 90235,
2	モバイル	39	28215, 35671, ...	
3	液晶	25	59321, 38106, ...	
4	プリンタ	14	67391, 984, ...	
:	:	:	:	
:	:	:	:	

GTBL 701 702 703 704 705

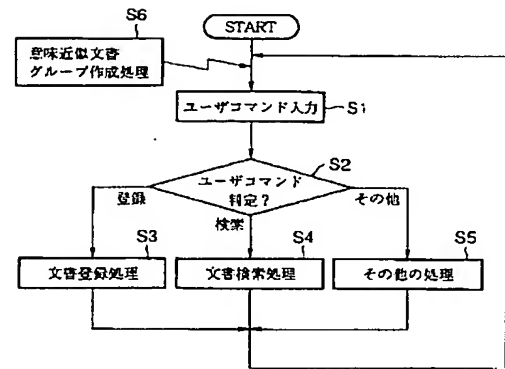
【図10】



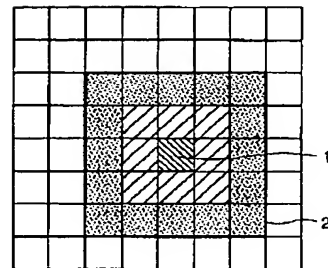
【図11】



【図9】



【図13】



【図14】

